

**This Page Is Inserted by IFW Operations  
and is not a part of the Official Record**

## **BEST AVAILABLE IMAGES**

**Defective images within this document are accurate representations of the original documents submitted by the applicant.**

**Defects in the images may include (but are not limited to):**

- **BLACK BORDERS**
- **TEXT CUT OFF AT TOP, BOTTOM OR SIDES**
- **FADED TEXT**
- **ILLEGIBLE TEXT**
- **SKEWED/SLANTED IMAGES**
- **COLORED PHOTOS**
- **BLACK OR VERY BLACK AND WHITE DARK PHOTOS**
- **GRAY SCALE DOCUMENTS**

**IMAGES ARE BEST AVAILABLE COPY.**

**As rescanning documents *will not* correct images,  
please do not report the images to the  
Image Problem Mailbox.**

**AUTOMATIC DOCUMENT CLASSIFICATION DEVICE, LEARNING DEVICE,  
CLASSIFICATION DEVICE, AUTOMATIC DOCUMENT CLASSIFICATION  
METHOD, LEARNING METHOD, CLASSIFICATION METHOD AND STORAGE  
MEDIUM**

Patent Number: JP11085796  
Publication date: 1999-03-30  
Inventor(s): OTANI NORIKO; ITO SHIRO; SHIBATA SHOGO; UEDA TAKANARI; IKEDA YUJI  
Applicant(s): CANON INC  
Requested Patent: ☐ JP11085796  
Application JP19970250125 19970901  
Priority Number(s):  
IPC Classification: G06F17/30 ; G06F17/21  
EC Classification:  
Equivalents:

**Abstract**

**PROBLEM TO BE SOLVED:** To provide an automatic document classification device which can appropriately classify a document where other topics different from a subject appear.

**SOLUTION:** The automatic document classification device refers to a valid word dictionary and obtains paragraph vectors on a learning document and a document being a classification object (paragraph vector calculation part 105). An other topic paragraph is decided from the distribution of the paragraph vectors (other topic paragraph decision part 107) and the valid paragraph vector is taken out from the paragraph vectors by referring to the other topic paragraph. A document vector is obtained from the paragraph vector (document vector calculation part 109). In a learning phase, the folder vectors of respective categories are obtained by using the document vector of the learning document (folder vector calculation part 111). In a classification phase, the category to which the document of the classification document belongs is decided in accordance with a comparison result between the document vector of the document being the classification object and the folder vectors of the respective categories (classification decision part 113).

Data supplied from the esp@cenet database - I2

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平11-85796

(43) 公開日 平成11年(1999) 3月30日

(51) Int.Cl.<sup>6</sup>

識別記号

F I

G 0 6 F 17/30

G 0 6 F 15/401

3 1 0 D

17/21

15/20

5 9 0 E

審査請求 未請求 請求項の数 9 F D (全 14 頁)

(21) 出願番号 特願平9-250125

(22) 出願日 平成9年(1997) 9月1日

(71) 出願人 000001007

キヤノン株式会社

東京都大田区下丸子3丁目30番2号

(72) 発明者 大谷 紀子

東京都大田区下丸子3丁目30番2号 キヤ  
ノン株式会社内

(72) 発明者 伊藤 史朗

東京都大田区下丸子3丁目30番2号 キヤ  
ノン株式会社内

(72) 発明者 柴田 昇吾

東京都大田区下丸子3丁目30番2号 キヤ  
ノン株式会社内

(74) 代理人 弁理士 渡部 敏彦

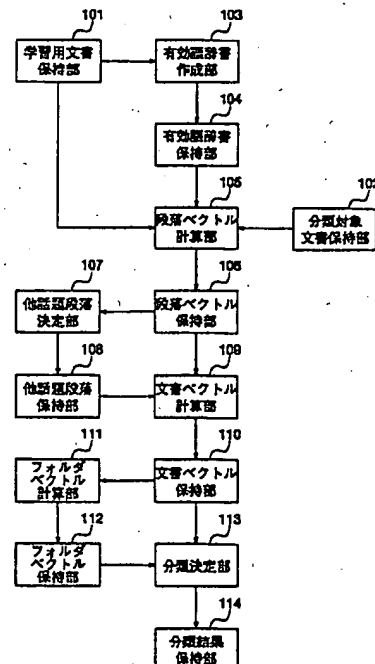
最終頁に続く

(54) 【発明の名称】 文書自動分類装置、学習装置、分類装置、文書自動分類方法、学習方法、分類方法および記憶媒体

(57) 【要約】

【課題】 主題と異なる他話題が出現する文書に対してその分類を適正に行うことができる文書自動分類装置を提供する。

【解決手段】 文書自動分類装置は、学習用文書と分類対象文書とのそれぞれについて、有効語辞書を参照して段落ベクトルを求め（段落ベクトル計算部105）、その段落ベクトルの分布から他話題段落を決定し（他話題段落決定部107）、その他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、この段落ベクトルから文書ベクトルを求める（文書ベクトル計算部109）。学習フェーズでは、学習用文書の文書ベクトルを用いて各カテゴリのフォルダベクトルを求め（フォルダベクトル計算部111）、分類フェーズでは、分類対象文書の文書ベクトルと各カテゴリのフォルダベクトルとの比較結果に応じて分類対象文書が属するカテゴリを決定する（分類決定部113）。



## 【特許請求の範囲】

【請求項1】 学習用文書と該学習用文書から選出された有効語を集めて作成した有効語辞書とを用いて、分類対象文書をユーザの意図に沿って分類する文書自動分類装置において、前記学習用文書と前記分類対象文書とのそれぞれについて、前記有効語辞書を参照して段落ベクトルを求める段落ベクトル計算手段と、前記学習用文書と前記分類対象文書とのそれぞれについて、その段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定する他話題段落決定手段と、前記学習用文書と前記分類対象文書とのそれぞれについて、その他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求める文書ベクトル計算手段と、前記学習用文書について求められた文書ベクトルを用いて各カテゴリのフォルダベクトルを求めるフォルダベクトル計算手段と、前記分類対象文書について求められた文書ベクトルと前記各カテゴリのフォルダベクトルとを比較し、該比較結果に応じて前記分類対象文書が属するカテゴリを決定する分類決定手段とを備えることを特徴とする文書自動分類装置。

【請求項2】 分類対象文書をユーザの意図に沿って分類する文書自動分類システムに用いられる、前記分類対象文書が属するカテゴリを決定するための基準を求めるための学習装置において、学習用文書を保持する学習用文書保持手段と、前記学習用文書から有効語を選出し、該選出された有効語を集めて有効語辞書を作成する有効語辞書作成手段と、前記学習用文書について前記有効語辞書を参照して段落ベクトルを求める段落ベクトル計算手段と、前記学習用文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定する他話題段落決定手段と、前記学習用文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求める文書ベクトル計算手段と、前記学習用文書の文書ベクトルを用いて前記分類対象文書が属するカテゴリを決定するための基準となる各カテゴリのフォルダベクトルを求めるフォルダベクトル計算手段とを備えることを特徴とする学習装置。

【請求項3】 分類対象文書をユーザの意図に沿って分類する文書自動分類システムに請求項2記載の学習装置とともに用いられる、前記分類対象文書が属するカテゴリを決定するための分類装置において、前記分類対象文書を保持する分類対象文書保持手段と、前記分類対象文書について前記有効語辞書を参照して段落ベクトルを求める段落ベクトル計算手段と、前記分類対象文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定する他話題段落決定手段と、前記分類対象文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトル

を用いて文書ベクトルを求める文書ベクトル計算手段と、前記分類対象文書の文書ベクトルと前記各カテゴリのフォルダベクトルとを比較し、該比較結果に応じて前記分類対象文書が属するカテゴリを決定する分類決定手段とを備えることを特徴とする分類装置。

【請求項4】 学習用文書と該学習用文書から選出された有効語を集めて作成した有効語辞書とを用いて、分類対象文書をユーザの意図に沿って分類する文書自動分類方法において、前記分類対象文書が属するカテゴリを決定するための基準を求めるための学習工程と、前記基準を用いて前記分類対象文書が属するカテゴリを決定するための分類工程とを有し、前記学習工程は、前記学習用文書について前記有効語辞書を参照して段落ベクトルを求める工程と、前記学習用文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定する工程と、前記学習用文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求める工程と、前記学習用文書の文書ベクトルを用いて前記分類対象文書が属するカテゴリを決定するための基準となる各カテゴリのフォルダベクトルを求める工程とを含み、前記分類工程は、前記分類対象文書について前記有効語辞書を参照して段落ベクトルを求める工程と、前記分類対象文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定する工程と、前記分類対象文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求める工程と、前記分類対象文書の文書ベクトルと前記各カテゴリのフォルダベクトルとを比較し、該比較結果に応じて前記分類対象文書が属するカテゴリを決定する工程とを含むことを特徴とする文書自動分類方法。

【請求項5】 分類対象文書をユーザの意図に沿って分類する文書自動分類システムに用いられる、前記分類対象文書が属するカテゴリを決定するための基準を求めるための学習方法において、学習用文書を保持する工程と、前記学習用文書から有効語を選出し、該選出された有効語を集めて有効語辞書を作成する工程と、前記学習用文書について前記有効語辞書を参照して段落ベクトルを求める工程と、前記学習用文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定する工程と、前記学習用文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求める工程と、前記学習用文書の文書ベクトルを用いて前記分類対象文書が属するカテゴリを決定するための基準となる各カテゴリのフォルダベクトルを求める工程とを含むことを特徴とする学習方法。

【請求項6】 分類対象文書をユーザの意図に沿って分類する文書自動分類システムに請求項5記載の学習方法

とともに用いられる、前記分類対象文書が属するカテゴリを決定するための分類方法において、前記分類対象文書を保持する工程と、前記分類対象文書について前記有効語辞書を参照して段落ベクトルを求める工程と、前記分類対象文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定する工程と、前記分類対象文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求める工程と、前記分類対象文書の文書ベクトルと前記各カテゴリのフォルダベクトルとを比較し、該比較結果に応じて前記分類対象文書が属するカテゴリを決定する工程とを含むことを特徴とする分類方法。

【請求項7】 学習用文書と該学習用文書から選出された有効語を集めて作成した有効語辞書とを用いて、分類対象文書をユーザの意図に沿って分類する文書自動分類システムを構築するためのプログラムを格納した記憶媒体において、前記プログラムは、前記分類対象文書が属するカテゴリを決定するための基準を求めるための学習プログラムと、前記基準を用いて前記分類対象文書が属するカテゴリを決定するための分類プログラムとを有し、前記学習プログラムは、前記学習用文書について前記有効語辞書を参照して段落ベクトルを求めるモジュールと、前記学習用文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定するモジュールと、前記学習用文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求めるモジュールと、前記学習用文書の文書ベクトルを用いて前記分類対象文書が属するカテゴリを決定するための基準となる各カテゴリのフォルダベクトルを求めるモジュールとを含み、前記分類プログラムは、前記分類対象文書について前記有効語辞書を参照して段落ベクトルを求めるモジュールと、前記分類対象文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定するモジュールと、前記分類対象文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求めるモジュールと、前記分類対象文書の文書ベクトルと前記各カテゴリのフォルダベクトルとを比較し、該比較結果に応じて前記分類対象文書が属するカテゴリを決定するモジュールとを含むことを特徴とする記憶媒体。

【請求項8】 分類対象文書をユーザの意図に沿って分類する文書自動分類システムに用いられ、前記分類対象文書が属するカテゴリを決定するための基準を求める学習装置を構築するための学習プログラムを格納した記憶媒体において、前記学習プログラムは、学習用文書を保持するモジュールと、前記学習用文書から有効語を選出し、該選出された有効語を集めて有効語辞書を作成するモジュールと、前記学習用文書について前記有効語辞書

を参照して段落ベクトルを求めるモジュールと、前記学習用文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定するモジュールと、前記学習用文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求めるモジュールと、前記学習用文書の文書ベクトルを用いて前記分類対象文書が属するカテゴリを決定するための基準となる各カテゴリのフォルダベクトルを求めるモジュールとを含むことを特徴とする記憶媒体。

【請求項9】 分類対象文書をユーザの意図に沿って分類する文書自動分類システムに請求項8記載の記憶媒体とともに用いられる、前記分類対象文書が属するカテゴリを決定する分類装置を構築するための分類プログラムを格納した記憶媒体において、前記分類プログラムは、前記分類対象文書を保持するモジュールと、前記分類対象文書について前記有効語辞書を参照して段落ベクトルを求めるモジュールと、前記分類対象文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定するモジュールと、前記分類対象文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求めるモジュールと、前記分類対象文書の文書ベクトルと前記各カテゴリのフォルダベクトルとを比較し、該比較結果に応じて前記分類対象文書が属するカテゴリを決定するモジュールとを含むことを特徴とする記憶媒体。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、分類対象文書をユーザの意図に沿って分類する文書自動分類装置、それに用いられる学習装置および分類装置と、文書自動分類方法、それに用いられる学習方法および分類方法と、文書自動分類装置を構築するための記憶媒体とに関する。

【0002】

【従来の技術】分類対象文書をユーザの意図に沿って分類する方法の一つとして、ベクトル空間モデルを利用した方法がある。このベクトル空間モデルでは、分類に有用な語や文書、カテゴリをベクトルで表現し、ベクトルの方向から文書が属するカテゴリを決定する。このベクトル空間モデルを利用した文書自動分類処理は、学習フェーズと分類フェーズとに分けられる。学習フェーズでは、予め正しく分類された学習用文書から分類に有用な語（以下、有効語という）を選出し、各有効語をベクトル表現して有効語辞書を作成する。また、学習用文書をベクトル表現して、各カテゴリの特徴を表すフォルダベクトルを算出する。分類フェーズでは、学習フェーズで得られた有効語辞書を用いて分類対象文書をベクトルで表現し（以下、文書ベクトルという）、この文書ベクトルとフォルダベクトルとを比較し、該比較結果に応じて

分類対象文書が属するカテゴリを決定する。

【0003】この方法を採用した文書自動分類装置の構成について図8ないし図11を参照しながら説明する。図8は従来の文書自動分類装置の構成を示すブロック図、図9は図8の文書自動分類装置における学習フェーズの処理手順を示すフローチャート、図10は図8の文書自動分類装置における分類フェーズの処理手順を示すフローチャート、図11は図8の文書自動分類装置における分類フェーズで求められた文書ベクトルの例を示す図である。

【0004】文書自動分類装置は、図8に示すように、学習用文書を保持する学習用文書保持部501と、分類対象文書を保持する分類対象文書保持部502と、学習用文書から有効語を選出し、この有効語を集めて有効語辞書を作成する有効語辞書作成部503と、有効語辞書を保持する有効語辞書保持部504と、学習用文書と分類対象文書とのそれぞれについて、有効語辞書を参照して文書ベクトルを求める文書ベクトル計算部505と、学習用文書と分類対象文書とのそれぞれについて求められた文書ベクトルを保持する文書ベクトル保持部506とを備える。

【0005】文書ベクトル保持部506に保持された学習用文書の文書ベクトルはフォルダベクトル計算部507に与えられ、フォルダベクトル計算部507は学習用文書の文書ベクトルを用いて各カテゴリのフォルダベクトルを求める。求められた各カテゴリのフォルダベクトルは、フォルダベクトル保持部508に保持される。

【0006】フォルダベクトル保持部508に保持された各カテゴリのフォルダベクトルは、文書ベクトル保持部506に保持された分類対象文書の文書ベクトルとともに分類決定部509に与えられ、分類決定部509は分類対象文書の文書ベクトルと各カテゴリのフォルダベクトルとを比較し、該比較結果に応じて分類対象文書が属するカテゴリを決定する。この決定された分類対象文書のカテゴリは分類結果保持部510に保持される。

【0007】次に、文書自動分類装置における学習フェーズの処理手順について図9を参照しながら説明する。

【0008】まず、ステップS601において学習用文書に含まれる語の中から、分類に有用な語を有効語として選定し、続くステップS602で、選定した有効語語が出現頻度や他の有効語との共起状況などによりベクトル表現し、有効語辞書として保持する。

【0009】次いで、ステップS603に進み、学習用文書から有効語を抽出し、続くステップS604で、有効語辞書を参照して取り出した有効語のベクトルの平均を取り、このベクトルの平均から学習用文書の文書ベクトルを求める。そして、ステップS605で、学習用文書における各カテゴリに属する文書の文書ベクトルの平均を取り、該文書のベクトルの平均からフォルダベクトルを求め、本処理を終了する。

【0010】この学習フェーズが終了すると、分類フェーズが開始される。この分類フェーズの処理手順について図10を参照しながら説明する。

【0011】分類フェーズでは、まずステップS701において分類対象文書から有効語を取り出し、続くステップS702で有効語辞書を参照して取り出した有効語のベクトルの平均を取り、このベクトルの平均から分類対象文書の文書ベクトルを求める。

【0012】次いで、ステップS703に進み、分類対象文書の文書ベクトルと学習フェーズで求められたフォルダベクトルとを比較し、該比較結果に応じて分類対象文書が属するカテゴリを決定し、本処理を終了する。

【0013】

【発明が解決しようとする課題】しかし、上述した従来の文書自動分類装置では、学習用文書または分類対象文書から有効語を取り出す際に、それぞれの文書中出现する全ての有効語を取り出し、それぞれの文書についてその全ての有効語を用いて文書ベクトルを求めるから、文書中に主題と異なる他話題が挿入されている場合には、文書ベクトルが主題からその方向を示すことがある。例えば、図11に示すように、有効語のベクトルa～fを有する文書においては、文書ベクトルが他の話題に出現する有効語のベクトルe、fに引っ張られ、主題に出現する有効語のベクトルa～dの方向（主題の方向）から文書ベクトルがそれと異なり、適正に文書の分類を行なうことができない。

【0014】本発明の目的は、主題と異なる他話題が出現する文書に対してその分類を適正に行うことができる文書自動分類装置、文書自動分類方法および記憶媒体を提供することにある。

【0015】本発明の他の目的は、主題と異なる他話題が出現する文書に対してその分類を適正に行うことが可能な文書自動分類システムを実現することができる学習装置、分類装置、学習方法、分類方法および記憶媒体を提供することにある。

【0016】

【課題を解決するための手段】請求項1記載の発明は、学習用文書と該学習用文書から選出された有効語を集めて作成した有効語辞書とを用いて、分類対象文書をユーザの意図に沿って分類する文書自動分類装置において、前記学習用文書と前記分類対象文書とのそれぞれについて、前記有効語辞書を参照して段落ベクトルを求める段落ベクトル計算手段と、前記学習用文書と前記分類対象文書とのそれぞれについて、その段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定する他話題段落決定手段と、前記学習用文書と前記分類対象文書とのそれぞれについて、その他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求める文書ベクトル計算手段と、前記学習用文書について求

められた文書ベクトルを用いて各カテゴリのフォルダベクトルを求めるフォルダベクトル計算手段と、前記分類対象文書について求められた文書ベクトルと前記各カテゴリのフォルダベクトルとを比較し、該比較結果に応じて前記分類対象文書が属するカテゴリを決定する分類決定手段とを備えることを特徴とする。

【0017】請求項2記載の発明は、分類対象文書をユーザの意図に沿って分類する文書自動分類システムに用いられる、前記分類対象文書が属するカテゴリを決定するための基準を求めるための学習装置において、学習用文書を保持する学習用文書保持手段と、前記学習用文書から有効語を選出し、該選出された有効語を集めて有効語辞書を作成する有効語辞書作成手段と、前記学習用文書について前記有効語辞書を参照して段落ベクトルを求める段落ベクトル計算手段と、前記学習用文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定する他話題段落決定手段と、前記学習用文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求める文書ベクトル計算手段と、前記学習用文書の文書ベクトルを用いて前記分類対象文書が属するカテゴリを決定するための基準となる各カテゴリのフォルダベクトルを求めるフォルダベクトル計算手段とを備えることを特徴とする。

【0018】請求項3記載の発明は、分類対象文書をユーザの意図に沿って分類する文書自動分類システムに請求項2記載の学習装置とともに用いられる、前記分類対象文書が属するカテゴリを決定するための分類装置において、前記分類対象文書を保持する分類対象文書保持手段と、前記分類対象文書について前記有効語辞書を参照して段落ベクトルを求める段落ベクトル計算手段と、前記分類対象文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定する他話題段落決定手段と、前記分類対象文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求める文書ベクトル計算手段と、前記分類対象文書の文書ベクトルと前記各カテゴリのフォルダベクトルとを比較し、該比較結果に応じて前記分類対象文書が属するカテゴリを決定する分類決定手段とを備えることを特徴とする。

【0019】請求項4記載の発明は、学習用文書と該学習用文書から選出された有効語を集めて作成した有効語辞書とを用いて、分類対象文書をユーザの意図に沿って分類する文書自動分類方法において、前記分類対象文書が属するカテゴリを決定するための基準を求めるための学習工程と、前記基準を用いて前記分類対象文書が属するカテゴリを決定するための分類工程とを有し、前記学習工程は、前記学習用文書について前記有効語辞書を参照して段落ベクトルを求める工程と、前記学習用文書の段落ベクトルの分布から主題とは異なる話題を表す他

話題段落を決定する工程と、前記学習用文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求める工程と、前記学習用文書の文書ベクトルを用いて前記分類対象文書が属するカテゴリを決定するための基準となる各カテゴリのフォルダベクトルを求める工程とを含み、前記分類工程は、前記分類対象文書について前記有効語辞書を参照して段落ベクトルを求める工程と、前記分類対象文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定する工程と、前記分類対象文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求める工程と、前記分類対象文書の文書ベクトルと前記各カテゴリのフォルダベクトルとを比較し、該比較結果に応じて前記分類対象文書が属するカテゴリを決定する工程とを含むことを特徴とする。

【0020】請求項5記載の発明は、分類対象文書をユーザの意図に沿って分類する文書自動分類システムに用いられる、前記分類対象文書が属するカテゴリを決定するための基準を求めるための学習方法において、学習用文書を保持する工程と、前記学習用文書から有効語を選出し、該選出された有効語を集めて有効語辞書を作成する工程と、前記学習用文書について前記有効語辞書を参照して段落ベクトルを求める工程と、前記学習用文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定する工程と、前記学習用文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求める工程と、前記学習用文書の文書ベクトルを用いて前記分類対象文書が属するカテゴリを決定するための基準となる各カテゴリのフォルダベクトルを求める工程とを含むことを特徴とする。

【0021】請求項6記載の発明は、分類対象文書をユーザの意図に沿って分類する文書自動分類システムに請求項5記載の学習方法とともに用いられる、前記分類対象文書が属するカテゴリを決定するための分類方法において、前記分類対象文書を保持する工程と、前記分類対象文書について前記有効語辞書を参照して段落ベクトルを求める工程と、前記分類対象文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定する工程と、前記分類対象文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求める工程と、前記分類対象文書の文書ベクトルと前記各カテゴリのフォルダベクトルとを比較し、該比較結果に応じて前記分類対象文書が属するカテゴリを決定する工程とを含むことを特徴とする。

【0022】請求項7記載の発明は、学習用文書と該学習用文書から選出された有効語を集めて作成した有効語

辞書とを用いて、分類対象文書をユーザの意図に沿って分類する文書自動分類システムを構築するためのプログラムを格納した記憶媒体において、前記プログラムは、前記分類対象文書が属するカテゴリを決定するための基準を求めるための学習プログラムと、前記基準を用いて前記分類対象文書が属するカテゴリを決定するための分類プログラムとを有し、前記学習プログラムは、前記学習用文書について前記有効語辞書を参照して段落ベクトルを求めるモジュールと、前記学習用文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定するモジュールと、前記学習用文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求めるモジュールと、前記学習用文書の文書ベクトルを用いて前記分類対象文書が属するカテゴリを決定するための基準となる各カテゴリのフォルダベクトルを求めるモジュールとを含み、前記分類プログラムは、前記分類対象文書について前記有効語辞書を参照して段落ベクトルを求めるモジュールと、前記分類対象文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定するモジュールと、前記分類対象文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求めるモジュールと、前記分類対象文書の文書ベクトルと前記各カテゴリのフォルダベクトルとを比較し、該比較結果に応じて前記分類対象文書が属するカテゴリを決定するモジュールとを含むことを特徴とする。

【0023】請求項8記載の発明は、分類対象文書をユーザの意図に沿って分類する文書自動分類システムに用いられ、前記分類対象文書が属するカテゴリを決定するための基準を求める学習装置を構築するための学習プログラムを格納した記憶媒体において、前記学習プログラムは、学習用文書を学習用文書保持手段に保持するモジュールと、前記学習用文書から有効語を選出し、該選出された有効語を集めて有効語辞書を作成するモジュールと、前記学習用文書について前記有効語辞書を参照して段落ベクトルを求めるモジュールと、前記学習用文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定するモジュールと、前記学習用文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求めるモジュールと、前記学習用文書の文書ベクトルを用いて前記分類対象文書が属するカテゴリを決定するための基準となる各カテゴリのフォルダベクトルを求めるモジュールとを含むことを特徴とする。

【0024】請求項9記載の発明は、分類対象文書をユーザの意図に沿って分類する文書自動分類システムに請求項8記載の記憶媒体とともに用いられる、前記分類対

象文書が属するカテゴリを決定する分類装置を構築するための分類プログラムを格納した記憶媒体において、前記分類プログラムは、前記分類対象文書を保持するモジュールと、前記分類対象文書について前記有効語辞書を参照して段落ベクトルを求めるモジュールと、前記分類対象文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定するモジュールと、前記分類対象文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求めるモジュールと、前記分類対象文書の文書ベクトルと前記各カテゴリのフォルダベクトルとを比較し、該比較結果に応じて前記分類対象文書が属するカテゴリを決定するモジュールとを含むことを特徴とする。

【0025】

【発明の実施の形態】以下に本発明の実施の形態について図を参照しながら説明する。

【0026】図1は本発明の文書自動分類装置の実施の一形態の機能構成を示すブロック図、図2は図1の文書自動分類装置のハードウェア構成を示すブロック図である。

【0027】文書自動分類装置は、図1に示すように、学習用文書を保持する学習用文書保持部101と、分類対象文書を保持する分類対象文書保持部102と、学習用文書から選出された有効語を集めて有効語辞書を作成する有効語辞書作成部103と、有効語辞書を保持する有効語辞書保持部104と、学習用文書と分類対象文書とのそれぞれについて、有効語辞書を参照して段落ベクトルを求める段落ベクトル計算部105と、学習用文書と分類対象文書とのそれぞれについて求められた段落ベクトルを保持する段落ベクトル保持部106とを備える。

【0028】段落ベクトル保持部106に保持された学習用文書と分類対象文書とのそれぞれの段落ベクトルは、他話題段落決定部107に与えられ、他話題段落決定部107は、学習用文書と分類対象文書とのそれぞれについて、その段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定する。この学習用文書と分類対象文書とのそれぞれについて決定された他話題段落は、他話題段落保持部108に保持される。

【0029】他話題段落保持部108に保持された学習用文書と分類対象文書とのそれぞれの他話題段落は、段落ベクトル保持部106に保持された学習用文書と分類対象文書とのそれぞれの段落ベクトルとともに文書ベクトル計算部109に与えられる。文書ベクトル計算部109は、学習用文書と分類対象文書とのそれぞれについて、その他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求める。学習用文書と分類対象文書とのそれぞれについて求められた文書ベクトル



ルは文書ベクトル保持部110に保持される。

【0030】文書ベクトル保持部110に保持された学習用文書の文書ベクトルはフォルダベクトル計算部111に与えられる。フォルダベクトル計算部111は学習用文書の文書ベクトルを用いて各カテゴリのフォルダベクトルを求め、求められた各カテゴリのフォルダベクトルはフォルダベクトル保持部112に保持される。

【0031】フォルダベクトル保持部112に保持された各カテゴリのフォルダベクトルは、文書ベクトル保持部110に保持された分類対象文書の文書ベクトルとともに分類決定部113に与えられる。分類決定部113は、分類対象文書の文書ベクトルと各カテゴリのフォルダベクトルとを比較し、該比較結果に応じて分類対象文書が属するカテゴリを決定し、この決定された分類対象文書のカテゴリは分類結果保持部114に保持される。

【0032】文書自動分類装置のハードウェア構成においては、図2に示すように、ROM201に格納されている制御プログラムを実行して後述する制御(図3および図4に示す制御)を行う中央処理装置203が設けられている。中央処理装置203の演算処理の作業領域としてはRAM202が用いられ、また、RAM202は、段落ベクトル保持部106、他話題段落保持部108、文書ベクトル保持部110、分類結果保持部114のための記憶領域を提供する。

【0033】中央処理装置203には、ROM201およびRAM202とともに、ハードディスク装置204がバス205を介して接続され、ハードディスク装置204は、学習用文書保持部101、分類対象文書保持部102、有効語辞書保持部104およびフォルダベクトル保持部112を構成する。なお、ハードディスク装置204に代えて、他の記憶媒体を用いて、学習用文書保持部101、分類対象文書保持部102、有効語辞書保持部104およびフォルダベクトル保持部112を構成することも可能である。

【0034】次に、本文書自動分類装置が実行する処理について図3および図4を参照しながら説明する。図3は図1の文書自動分類装置における学習フェーズの処理手順を示すフローチャート、図4は図1の文書自動分類装置における分類フェーズの処理手順を示すフローチャート、図5は図1の文書自動分類装置における文書ベクトルの例を示す図である。

【0035】本文書自動分類装置における処理は学習フェーズと分類フェーズとに分けられ、最初に、学習フェーズの処理手順について図3を参照しながら説明する。

【0036】学習フェーズでは、図3に示すように、まずステップS301において学習用文書に含まれる語の中から、分類に有用な有効語として選定し、続くステップS302で、選定した有効語を出現頻度や他の有効語との共起状況などによりベクトル表現し、有効語辞書として保持する。

【0037】次いで、ステップS303に進み、学習用文書から有効語を抽出し、続くステップS304で、有効語辞書を参照して取り出した有効語のベクトルの平均を段落毎に算出し、各段落の段落ベクトルとする。

【0038】次いで、ステップS305に進み、各文書毎に段落ベクトルの分布を調べ、他の段落ベクトルと極端に違う方向の段落ベクトルがあれば、その段落ベクトルの段落を主題とは異なる話題を表す他話題段落として決定する。また、各段落ベクトルにおいてそれぞれの向きが近接しているときには、全ての段落ベクトルが主題を表していると判断する。すなわち全ての段落を他話題段落ではないとする。例えば、各段落ベクトル毎に他の段落ベクトルとの余弦値の総和を求め、この総和が他より極端に小さい段落ベクトルがあれば、該段落ベクトルの段落を他話題段落として決定する。ここで、他の段落ベクトルとの余弦値の総和が正規分布に従うとすれば、該分布の数%以下を示す段落ベクトルの段落を他話題段落として定義することが可能である。

【0039】他話題段落を決定すると、ステップS306に進み、他話題段落の段落ベクトルを段落ベクトル群から除去し、残りの段落ベクトルの平均を取って文書ベクトルとする。続くステップS307では、学習用文書における各カテゴリに属する文書の文書ベクトルの平均を取り、該文書ベクトルの平均からフォルダベクトルを求め、本処理を終了する。

【0040】この学習フェーズが終了すると、分類フェーズが開始される。この分類フェーズの処理手順について図4を参照しながら説明する。

【0041】分類フェーズでは、図4に示すように、まずステップS401において分類対象文書から有効語を取り出し、続くステップS402で有効語辞書を参照して取り出した有効語のベクトルの平均を段落毎に取り、各段落のベクトルの平均から分類対象文書の文書ベクトルを求める。

【0042】次いで、ステップS403に進み、各文書毎に段落ベクトルの分布を調べて分類対象文書の他話題段落を決定する。この他話題段落の決定処理は上述の学習フェーズにおけるステップS305の処理内容と同じであり、その説明は省略する。分類対象文書の文書ベクトルと学習フェーズで求められたフォルダベクトルとを比較し、該比較結果に応じて分類対象文書が属するカテゴリを決定し、本処理を終了する。

【0043】例えば、図11に示すように、主題と異なる他話題が出現する文書において、有効語のベクトルa、bが段落Aに、有効語のベクトルc、dが段落Bに、有効語のベクトルe、fが段落Cにそれぞれ出現したとすると、各段落A、B、Cの段落ベクトルの内の段落Cが、図5に示すように、主題と異なる他話題を表す他話題段落として決定され、この段落Cの段落ベクトルを除去して各段落A、Bの段落ベクトルから文書ベクトル

ルが求められる。よって、求められた文書ベクトルが主題の方向にほぼ一致することになり、適正に文書の分類を行うことができる。

【0044】以上より、本実施の形態では、各段落の段落ベクトルを求めてその分布からはずれている段落を他話題段落として除去し、残りの段落から文書ベクトルを求めることにより、主題の方向をほぼ示すような文書ベクトルが得られ、主題と異なる他話題が出現する文書に対してその分類を適正に行うことができる。

【0045】なお、本実施の形態では、他話題段落を段落ベクトル間の余弦値の総和に基づき決定する例を示したが、これに限定されるものではなく、他の値を用いて他話題段落を決定することも可能である。また、段落ベクトルの分布において、該分布からはずれているか否かの決定に正規分布に従う値を基準値として決定する方法を示したが、これに限定されるものではない。

【0046】また、本実施の形態では、学習フェーズにおいて、段落ベクトルの計算、他話題段落の決定、文書ベクトルの計算の各処理を学習用文書に含まれる全ての文書に対して実行し、その後に次の処理を実行するように設定しているが、これに限定されるものではなく、1文書づつ各処理を実行するように設定することも可能である。

【0047】さらに、本実施の形態では、段落単位で話題を取り扱っているが、これに限定されるものではなく、文や節など、他の文章単位で扱うことも可能である。

【0048】さらに、本実施の形態では、文書ベクトルを平均してフォルダベクトルを求めるように説明しているが、各カテゴリの段落ベクトルの平均をフォルダベクトルとしてもよい。この場合、学習フェーズにおいて、文書ベクトルを求める必要はない。

【0049】さらに、本実施の形態では、上述の処理（各ブロックの機能）を実行するためのプログラムをROMに格納した例を示したが、他の記憶媒体を用いて上記プログラムを供給するように構成することも可能である。また、各ブロックの機能をそれぞれ有する回路構成により本装置を構成することも可能である。

【0050】さらに、本装置をコンピュータなどの情報処理装置上に構築することも可能である。この場合、上述の処理（各ブロックの機能）を実行するためのプログラムを格納した記憶媒体を準備し、CPUなどが該記憶媒体から上記プログラムを読み出して実行することにより、文書自動分類装置が構成される。上記プログラムを供給するための記憶媒体としては、フロッピーディスク、ハードディスク、光ディスク、光磁気ディスク、CDROM、CD-R、磁気テープ、不揮発性メモ리카ード、ROMなどを用いることができる。なお、上記プログラムの実行により文書自動分類装置を構成する場合には、コンピュータ上で稼働しているOSが上記プログラ

ムに含まれる処理の一部または全てを実行するように構成されている場合も含まれる。また、記憶媒体から供給されたプログラムがコンピュータに搭載された拡張機能ボードまたは接続された周辺拡張ユニットに書き込まれた後に、拡張機能ボードまたは周辺拡張ユニットに設けられたCPUが書き込まれたプログラムを実行する場合も含まれる。

【0051】さらに、本発明の原理は、複数の機器からなるシステム、ひとつの機器からなる装置のいずれにも適用することが可能である。

【0052】さらに、本実施の形態では、学習フェーズと分類フェーズとを一つの装置上で行う例を説明したが、これに限定されるものではなく、例えば、学習フェーズを行う装置と、分類フェーズを行う装置とを準備し、それぞれの装置を用いて文書の分類を行うように構成することもできる。この場合、学習フェーズを行う装置により、有効語辞書を作成したフォルダベクトルを求め、この有効語辞書およびフォルダベクトルを可搬記憶媒体または通信により、分類フェーズを行う装置に供給して分類を行う方法が用いられる。

【0053】この学習フェーズを行う装置および分類フェーズを行う装置について図6および図7を参照しながら説明する。図6は本発明の学習装置の実施の一形態の構成を示すブロック図、図7は本発明の分類装置の実施の一形態の構成を示すブロック図である。

【0054】学習フェーズを行う装置は、図6に示すように、学習用文書を保持する学習用文書保持部1001と、学習用文書から選出された有効語を集めて有効語辞書を作成する有効語辞書作成部1002と、有効語辞書を保持する有効語辞書保持部1003と、学習用文書について、有効語辞書を参照して段落ベクトルを求める段落ベクトル計算部1004と、学習用文書について求められた段落ベクトルを保持する段落ベクトル保持部1005とを備える。

【0055】段落ベクトル保持部1005に保持された学習用文書の段落ベクトルは他話題段落決定部1006に与えられ、他話題段落決定部1006は、学習用文書について、その段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定する。この学習用文書と分類対象文書とのそれぞれについて決定された他話題段落は、他話題段落保持部1007に保持される。

【0056】他話題段落保持部1007に保持された学習用文書の他話題段落は、段落ベクトル保持部1005に保持された学習用文書の段落ベクトルとともに文書ベクトル計算部1008に与えられる。文書ベクトル計算部1008は、学習用文書について、その他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求める。学習用文書について求められた文書ベクトルは文書ベクトル保持部1009に保持される。

【0057】文書ベクトル保持部1009に保持された学習用文書の文書ベクトルはフォルダベクトル計算部1010に与えられる。フォルダベクトル計算部1010は学習用文書の文書ベクトルを用いて各カテゴリのフォルダベクトルを求め、求められた各カテゴリのフォルダベクトルはフォルダベクトル保持部1011に保持される。

【0058】フォルダベクトル保持部1011に保持された各カテゴリのフォルダベクトル、および有効語辞書保持部1003に保持された有効語辞書は、可搬記憶媒体に記憶されて分類フェーズを行う装置に供給され、または通信により分類フェーズを行う装置に供給される。

【0059】分類フェーズを行う装置は、図7に示すように、分類対象文書を保持する分類対象文書保持部1101と、学習フェーズを行う装置から可搬記憶媒体または通信を介して供給された有効語辞書を保持する有効語辞書保持部1102と、学習フェーズを行う装置から可搬記憶媒体または通信を介して供給されたフォルダベクトルを保持するフォルダベクトル保持部1109と、分類対象文書について、有効語辞書を参照して段落ベクトルを求める段落ベクトル計算部1103と、分類対象文書について求められた段落ベクトルを保持する段落ベクトル保持部1104とを備える。

【0060】段落ベクトル保持部1104に保持された分類対象文書の段落ベクトルは他話題段落決定部1105に与えられ、他話題段落決定部1105は、分類対象文書について、その段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定する。この分類対象文書について決定された他話題段落は、他話題段落保持部1106に保持される。

【0061】他話題段落保持部1106に保持された分類対象文書の他話題段落は、段落ベクトル保持部1104に保持された分類対象文書の段落ベクトルとともに文書ベクトル計算部1107に与えられる。文書ベクトル計算部1107は、分類対象文書について、その他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求める。分類対象文書について求められた文書ベクトルは文書ベクトル保持部1108に保持される。

【0062】文書ベクトル保持部1108に保持された分類対象文書の文書ベクトルは、フォルダベクトル保持部1109に保持された各カテゴリのフォルダベクトルとともに分類決定部1110に与えられる。分類決定部1110は、分類対象文書の文書ベクトルと各カテゴリのフォルダベクトルとを比較し、該比較結果に応じて分類対象文書が属するカテゴリを決定し、この決定された分類対象文書のカテゴリは分類結果保持部1111に保持される。

【0063】

【発明の効果】以上に説明したように、請求項1記載の文書自動分類装置によれば、学習用文書と分類対象文書とのそれぞれについて、有効語辞書を参照して段落ベクトルを求める段落ベクトル計算手段と、学習用文書と分類対象文書とのそれぞれについて、その段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定する他話題段落決定手段と、学習用文書と分類対象文書とのそれぞれについて、その他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求める文書ベクトル計算手段と、学習用文書について求められた文書ベクトルを用いて各カテゴリのフォルダベクトルを求めるフォルダベクトル計算手段と、分類対象文書について求められた文書ベクトルと各カテゴリのフォルダベクトルとを比較し、該比較結果に応じて分類対象文書が属するカテゴリを決定する分類決定手段とを備えるから、主題の方向をほぼ示すような分類対象文書の文書ベクトルと、主題の方向を適正に示す各カテゴリのフォルダベクトルとが得られ、この文書ベクトルと各カテゴリのフォルダベクトルとの比較により、主題と異なる他話題が出現する文書に対してその分類を適正に行うことができる。

【0064】請求項2記載の学習装置によれば、学習用文書を保持する学習用文書保持手段と、学習用文書から有効語を選出し、該選出された有効語を集めて有効語辞書を作成する有効語辞書作成手段と、学習用文書について有効語辞書を参照して段落ベクトルを求める段落ベクトル計算手段と、学習用文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定する他話題段落決定手段と、学習用文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求める文書ベクトル計算手段と、学習用文書の文書ベクトルを用いて分類対象文書が属するカテゴリを決定するための基準となる各カテゴリのフォルダベクトルを求めるフォルダベクトル計算手段とを備えるから、主題の方向を適正に示す各カテゴリのフォルダベクトルが得られ、主題と異なる他話題が出現する文書に対してその分類を適正に行うことが可能な文書自動分類システムを実現することができる。

【0065】請求項3記載の分類装置によれば、分類対象文書を保持する分類対象文書保持手段と、分類対象文書について有効語辞書を参照して段落ベクトルを求める段落ベクトル計算手段と、分類対象文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定する他話題段落決定手段と、分類対象文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求める文書ベクトル計算手段と、分類対象文書の文書ベクトルと各カテゴリのフォルダベクトルとを

比較し、該比較結果に応じて分類対象文書が属するカテゴリを決定する分類決定手段とを備えるから、主題の方向をほぼ示すよう分類対象文書の文書ベクトルが得られ、この文書ベクトルと主題の方向を適正に示す各カテゴリのフォルダベクトルとの比較により、主題と異なる他話題が出現する文書に対してその分類を適正に行うことが可能な文書自動分類システムを実現することができる。

【0066】請求項4記載の文書自動分類方法によれば、分類対象文書が属するカテゴリを決定するための基準を求めるための学習工程と、基準を用いて分類対象文書が属するカテゴリを決定するための分類工程とを有し、学習工程が、学習用文書について有効語辞書を参照して段落ベクトルを求める工程と、学習用文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定する工程と、学習用文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求める工程と、学習用文書の文書ベクトルを用いて分類対象文書が属するカテゴリを決定するための基準となる各カテゴリのフォルダベクトルを求める工程とを含み、分類工程が、分類対象文書について有効語辞書を参照して段落ベクトルを求める工程と、分類対象文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定する工程と、分類対象文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求める工程と、分類対象文書の文書ベクトルと各カテゴリのフォルダベクトルとを比較し、該比較結果に応じて分類対象文書が属するカテゴリを決定する工程とを含むから、主題の方向をほぼ示すような分類対象文書の文書ベクトルと、主題の方向を適正に示す各カテゴリのフォルダベクトルとが得られ、この文書ベクトルと各カテゴリのフォルダベクトルとの比較により、主題と異なる他話題が出現する文書に対してその分類を適正に行うことができる。

【0067】請求項5記載の学習方法によれば、学習用文書を学習用文書保持手段に保持する工程と、学習用文書から有効語を選出し、該選出された有効語を集めて有効語辞書を作成する工程と、学習用文書について有効語辞書を参照して段落ベクトルを求める工程と、学習用文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定する工程と、学習用文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求める工程と、学習用文書の文書ベクトルを用いて分類対象文書が属するカテゴリを決定するための基準となる各カテゴリのフォルダベクトルを求める工程とを含むから、主題の方向を適正に示す各カテゴリのフォルダベクトルが得られ、主題と異なる他話題が出現す

る文書に対してその分類を適正に行うことが可能な文書自動分類システムを実現することができる。

【0068】請求項6記載の分類方法によれば、分類対象文書を保持する工程と、分類対象文書について有効語辞書を参照して段落ベクトルを求める工程と、分類対象文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定する工程と、分類対象文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求める工程と、分類対象文書の文書ベクトルと各カテゴリのフォルダベクトルとを比較し、該比較結果に応じて分類対象文書が属するカテゴリを決定する工程とを含むから、主題の方向をほぼ示すよう分類対象文書の文書ベクトルが得られ、この文書ベクトルと主題の方向を適正に示す各カテゴリのフォルダベクトルとの比較により、主題と異なる他話題が出現する文書に対してその分類を適正に行うことが可能な文書自動分類システムを実現することができる。

【0069】請求項7記載の記憶媒体によれば、前記プログラムが、分類対象文書が属するカテゴリを決定するための基準を求めるための学習プログラムと、基準を用いて分類対象文書が属するカテゴリを決定するための分類プログラムとを有し、学習プログラムが、学習用文書について有効語辞書を参照して段落ベクトルを求めるモジュールと、学習用文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定するモジュールと、学習用文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求めるモジュールと、学習用文書の文書ベクトルを用いて分類対象文書が属するカテゴリを決定するための基準となる各カテゴリのフォルダベクトルを求めるモジュールとを含み、分類プログラムが、分類対象文書について有効語辞書を参照して段落ベクトルを求めるモジュールと、分類対象文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定するモジュールと、分類対象文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求めるモジュールと、分類対象文書の文書ベクトルと各カテゴリのフォルダベクトルとを比較し、該比較結果に応じて前記分類対象文書が属するカテゴリを決定するモジュールとを含むから、主題の方向をほぼ示すような分類対象文書の文書ベクトルと、主題の方向を適正に示す各カテゴリのフォルダベクトルとが得られ、この文書ベクトルと各カテゴリのフォルダベクトルとの比較により、主題と異なる他話題が出現する文書に対してその分類を適正に行うことができる。

【0070】請求項8記載の記憶媒体によれば、学習プログラムが、学習用文書を保持するモジュールと、学習用文書から有効語を選出し、該選出された有効語を集め

て有効語辞書を作成するモジュールと、学習用文書について有効語辞書を参照して段落ベクトルを求めるモジュールと、学習用文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定するモジュールと、学習用文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求めるモジュールと、学習用文書の文書ベクトルを用いて分類対象文書が属するカテゴリを決定するための基準となる各カテゴリのフォルダベクトルを求めるモジュールとを含むから、主題の方向を適正に示す各カテゴリのフォルダベクトルが得られ、主題と異なる他話題が出現する文書に対してその分類を適正に行うことが可能な文書自動分類システムを実現することができる。

【0071】請求項9記載の記憶媒体によれば、分類プログラムが、分類対象文書を保持するモジュールと、分類対象文書について有効語辞書を参照して段落ベクトルを求めるモジュールと、分類対象文書の段落ベクトルの分布から主題とは異なる話題を表す他話題段落を決定するモジュールと、分類対象文書の他話題段落を参照してその段落ベクトルの中から有効な段落ベクトルを取り出し、該取り出した段落ベクトルを用いて文書ベクトルを求めるモジュールと、分類対象文書の文書ベクトルと前記各カテゴリのフォルダベクトルとを比較し、該比較結果に応じて前記分類対象文書が属するカテゴリを決定するモジュールとを含むから、主題の方向をほぼ示すよう分類対象文書の文書ベクトルが得られ、この文書ベクトルと主題の方向を適正に示す各カテゴリのフォルダベクトルとの比較により、主題と異なる他話題が出現する文書に対してその分類を適正に行うことが可能な文書自動分類システムを実現することができる。

【図面の簡単な説明】

【図1】本発明の文書自動分類装置の実施の一形態の機能構成を示すブロック図である。

【図2】図1の文書自動分類装置のハードウェア構成を示すブロック図である。

【図3】図1の文書自動分類装置における学習フェーズ

の処理手順を示すフローチャートである。

【図4】図1の文書自動分類装置における分類フェーズの処理手順を示すフローチャートである。

【図5】図1の文書自動分類装置における文書ベクトルの例を示す図である。

【図6】本発明の学習装置の実施の一形態の構成を示すブロック図である。

【図7】本発明の分類装置の実施の一形態の構成を示すブロック図である。

【図8】従来の文書自動分類装置の構成を示すブロック図である。

【図9】図8の文書自動分類装置における学習フェーズの処理手順を示すフローチャートである。

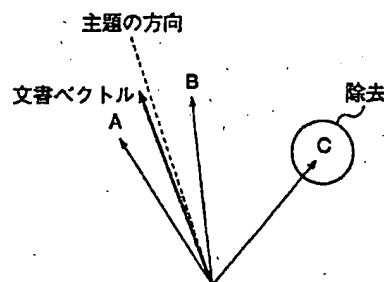
【図10】図8の文書自動分類装置における分類フェーズの処理手順を示すフローチャートである。

【図11】図8の文書自動分類装置における分類フェーズで求められた文書ベクトルの例を示す図である。

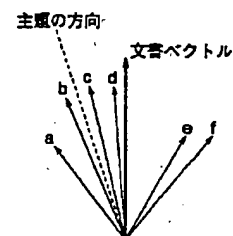
【符号の説明】

- 101, 1001 学習用文書保持部
- 102, 1101 分類対象文書保持部
- 103, 1002 有効語辞書作成部
- 104, 1003, 1102 有効語辞書保持部
- 105, 1004, 1103 段落ベクトル計算部
- 106, 1005, 1104 段落ベクトル保持部
- 107, 1006, 1105 他話題段落決定部
- 108, 1007, 1106 他話題段落保持部
- 109, 1008, 1107 文書ベクトル計算部
- 110, 1009, 1108 文書ベクトル保持部
- 111, 1010, フォルダベクトル計算部
- 112, 1011, 1109 フォルダベクトル保持部
- 113, 1110 分類決定部
- 114, 1111 分類結果保持部
- 201 ROM
- 202 RAM
- 203 中央処理装置
- 204 ハードディスク装置

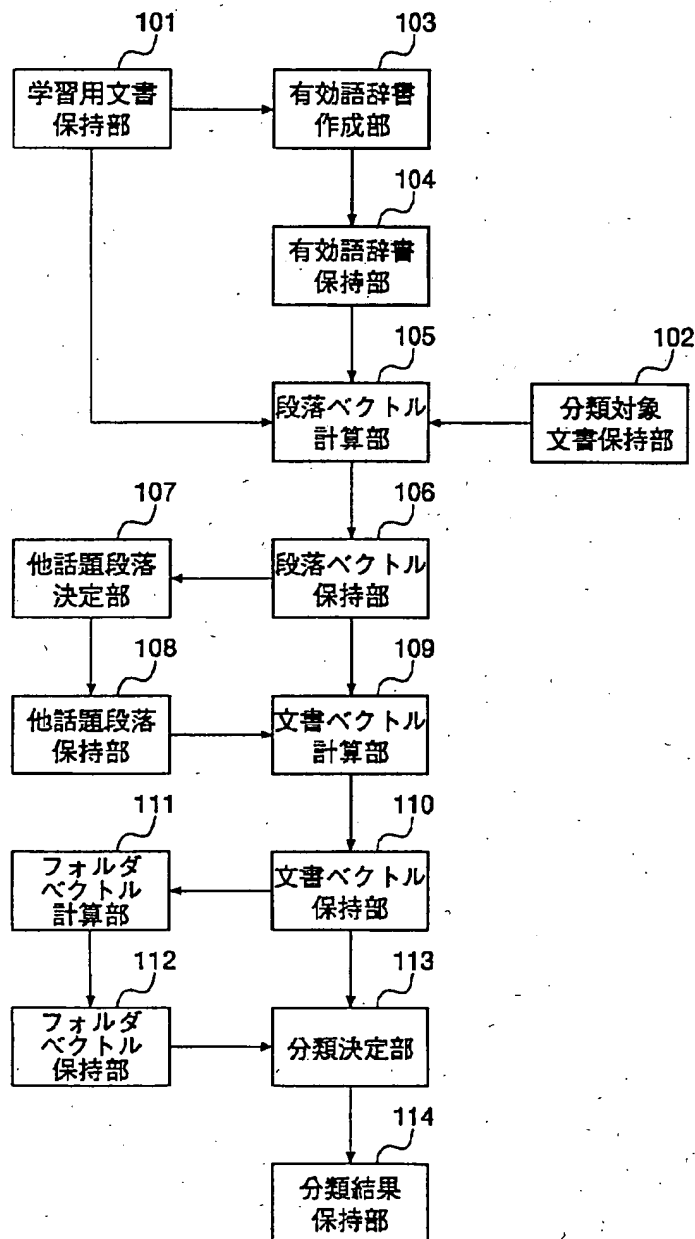
【図5】



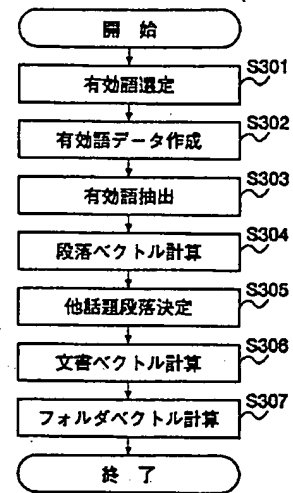
【図11】



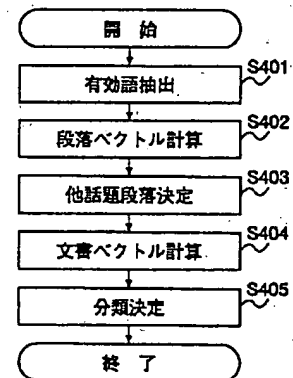
【図1】



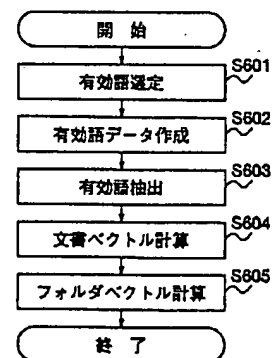
【図3】



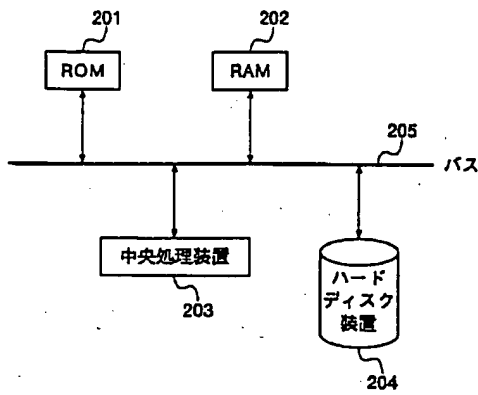
【図4】



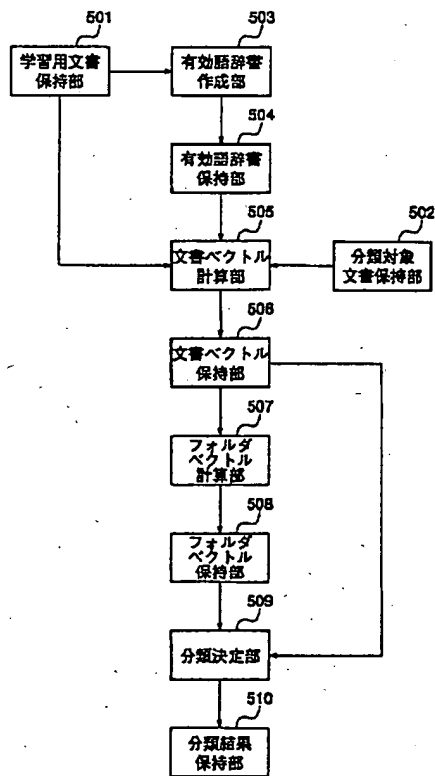
【図9】



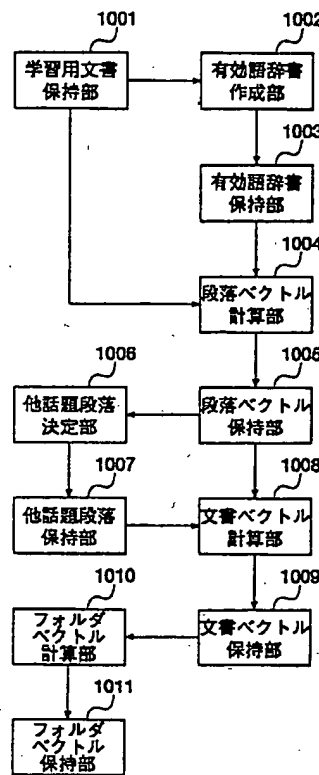
【図2】



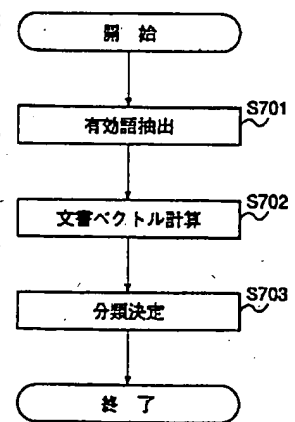
【図8】



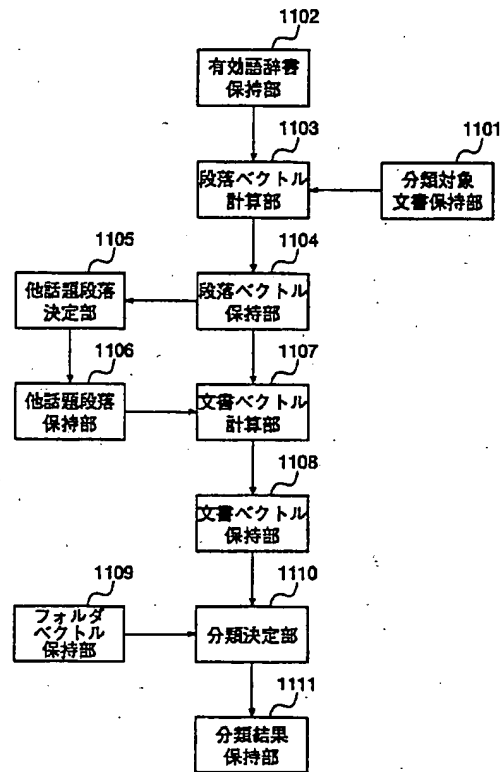
【図6】



【図10】



【図7】




---

フロントページの続き

(72)発明者 上田 隆也  
 東京都大田区下丸子3丁目30番2号 キヤ  
 ノン株式会社内

(72)発明者 池田 裕治  
 東京都大田区下丸子3丁目30番2号 キヤ  
 ノン株式会社内